**UNITED STATES DISTRICT COURT**
**SOUTHERN DISTRICT OF NEW YORK**

| | |
|---|---|
| IN RE:<br><br>OPENAI, INC.,<br>COPYRIGHT INFRINGEMENT LITIGATION<br><br>This Document Relates To:<br><br>All Actions | 25-md-3143 (SHS) (OTW)<br><br>Hon. Sidney H. Stein<br>Hon. Ona T. Wang |

<u>**NOTICE OF OPENAI'S TECHNOLOGY TUTORIAL DEMONSTRATIVES**</u>

# OpenAI's
# Technology Tutorial
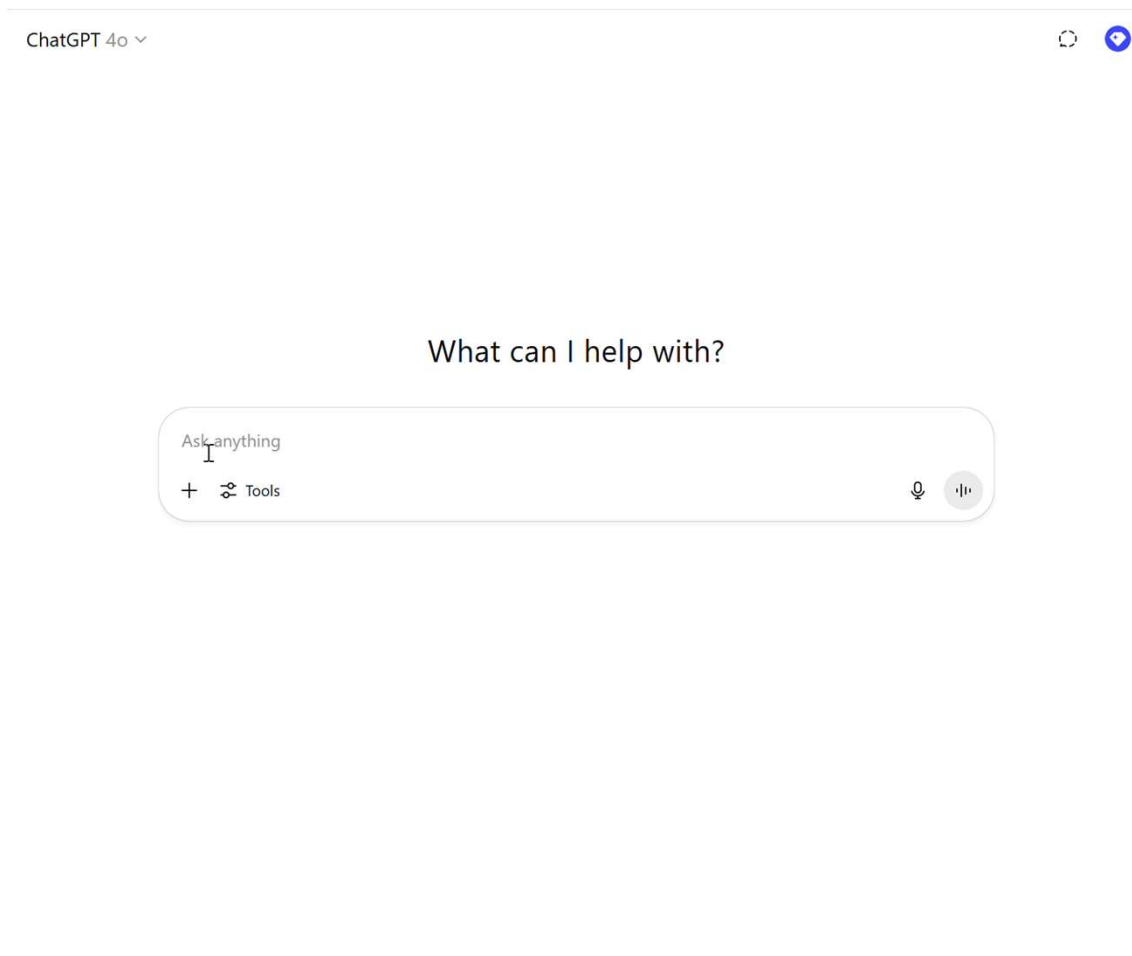
In re OpenAI Copyright MDL
MDL No. 25-md-3143 (S.D.N.Y.)

June 26, 2025

# What is a GPT?

**G**ENERATIVE

**P**RE-TRAINED

**T**RANSFORMER

2

ChatGPT 4o ∨

What can I help with?

Ask anything

+    Tools

# ChatGPT generates innovation

Through chatbots on their phones, farmers received customised advice on their chilli crops

NASA-GPT: Searching the Entire NASA Technical Reports Server Using AI

SUPERCOMPUTING

"The future applications within NASA are nearly endless."

XB-70 Image credit: US Air Force

UO

What is the size of the inlet bleed holes on the XB-70?

https://www.linkedin.com/posts/world-economic-forum_digital-agriculture-farmers-ugcPost-7279073700228808704-CMqR?utm_source=share&utm_medium=member_ios

4

# ChatGPT uses Large Language Models (LLMs): "GPTs"

A GPT is a statistical model that predicts the next word by learning patterns in language

**Prompt:**

How do humans learn language?

**Response:**

Through patterns, observation, and social

————

| | |
|---|---|
| Interactions | 23% |
| Feedback | 11% |
| Engagement | 10% |
| Media | 6% |

5

# What is a GPT?

**G**ENERATIVE — *Not a database; generates **new** content*
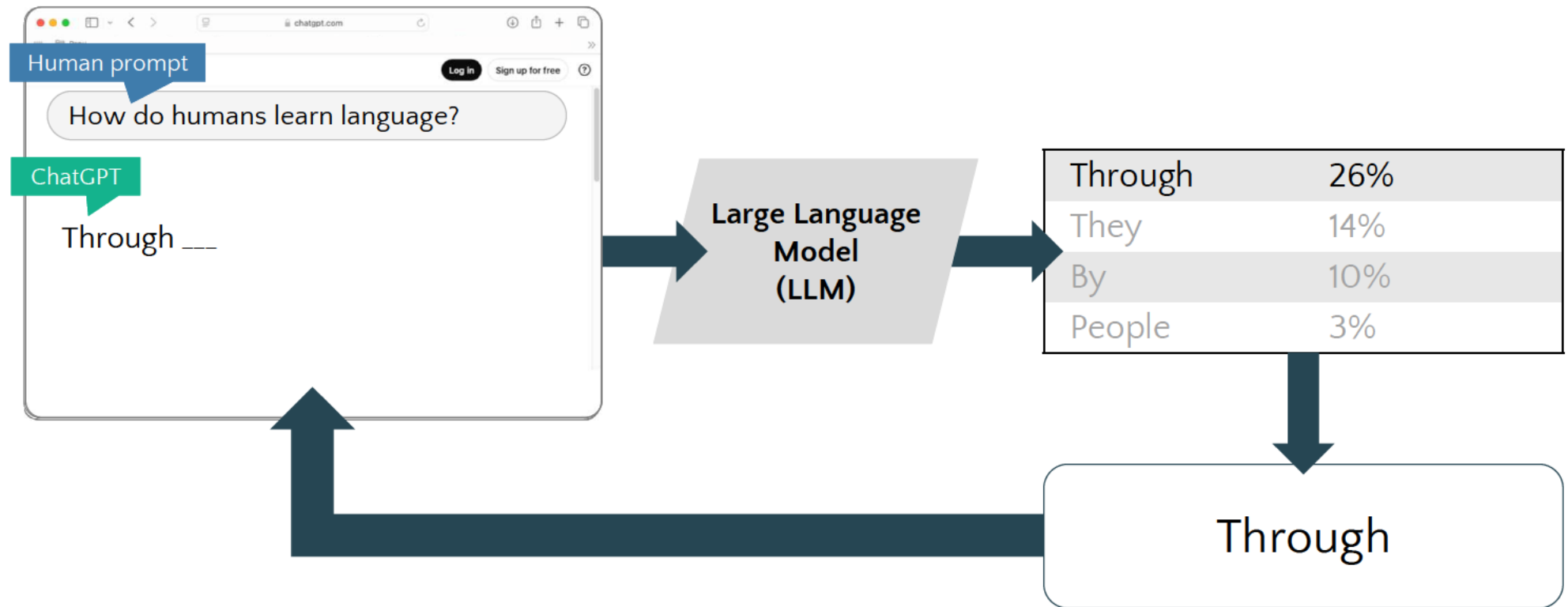
**P**RE-TRAINED — *Trained to identify linguistic patterns*

**T**RANSFORMER — *Transforms language data into a platform for **innovation***

6

# LLMs are not databases or libraries

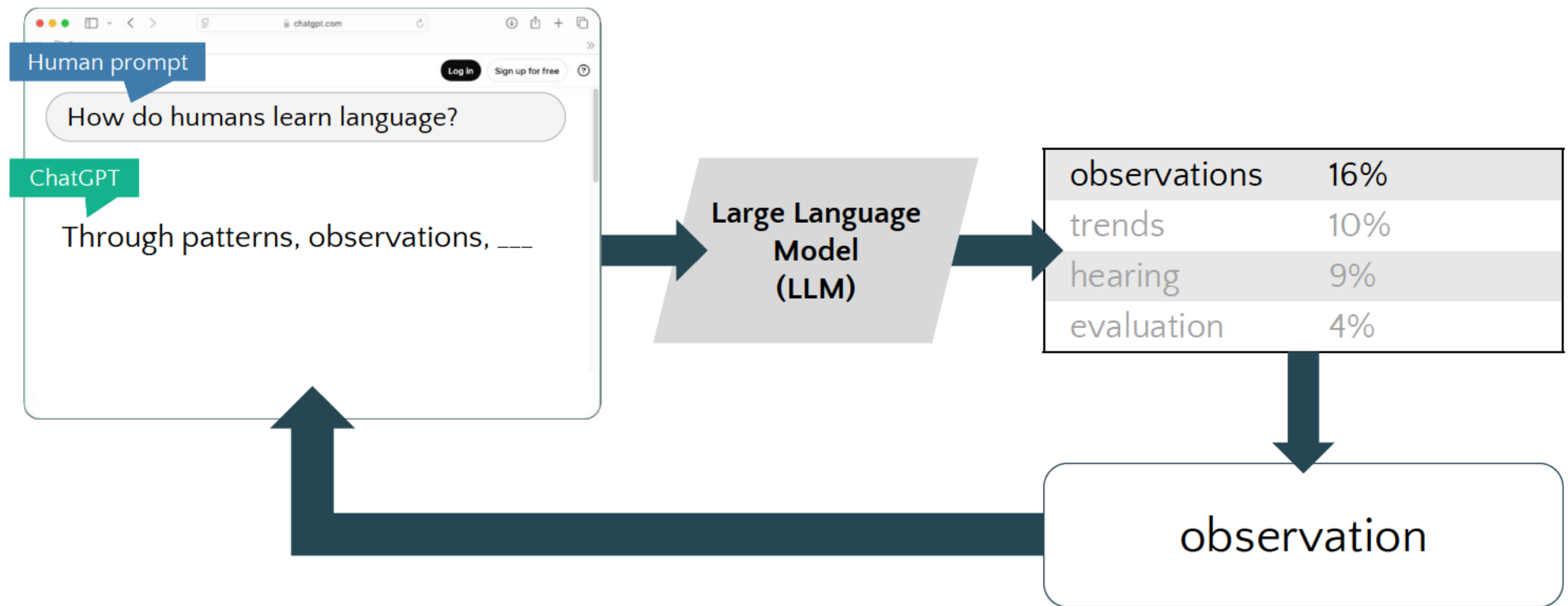They do not store information or retrieve information from a database



7

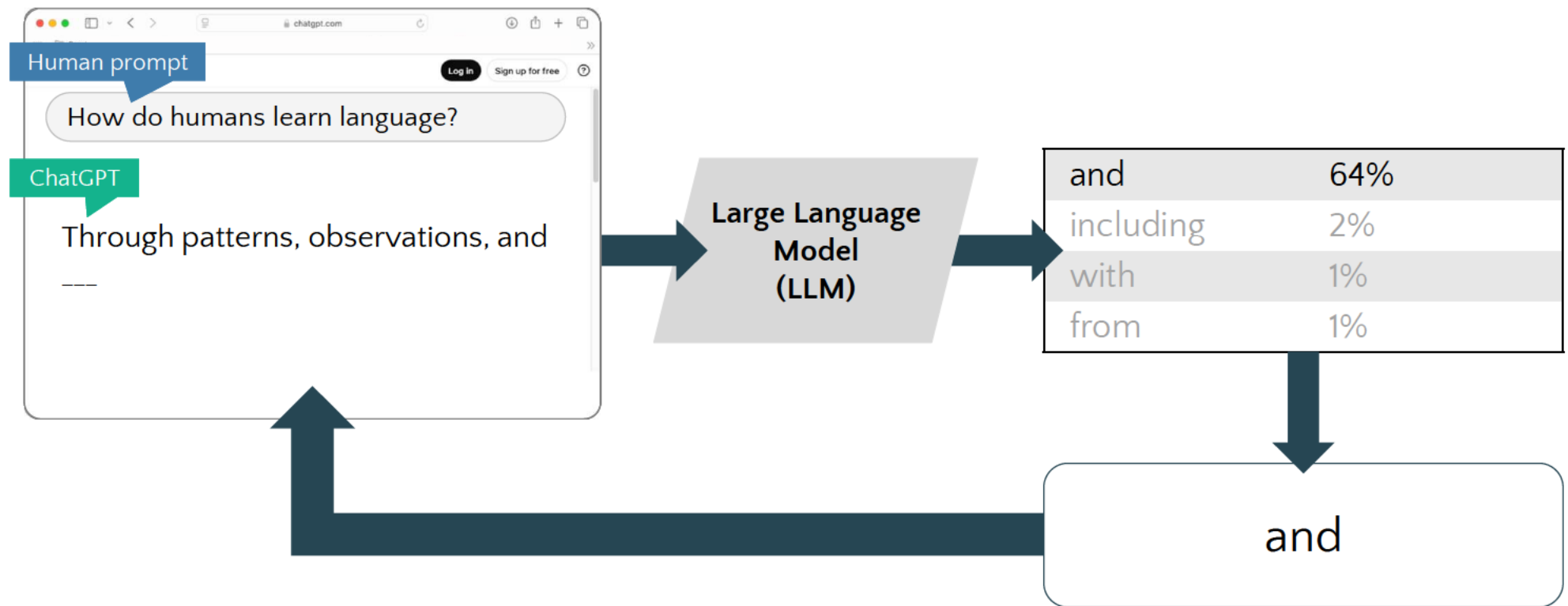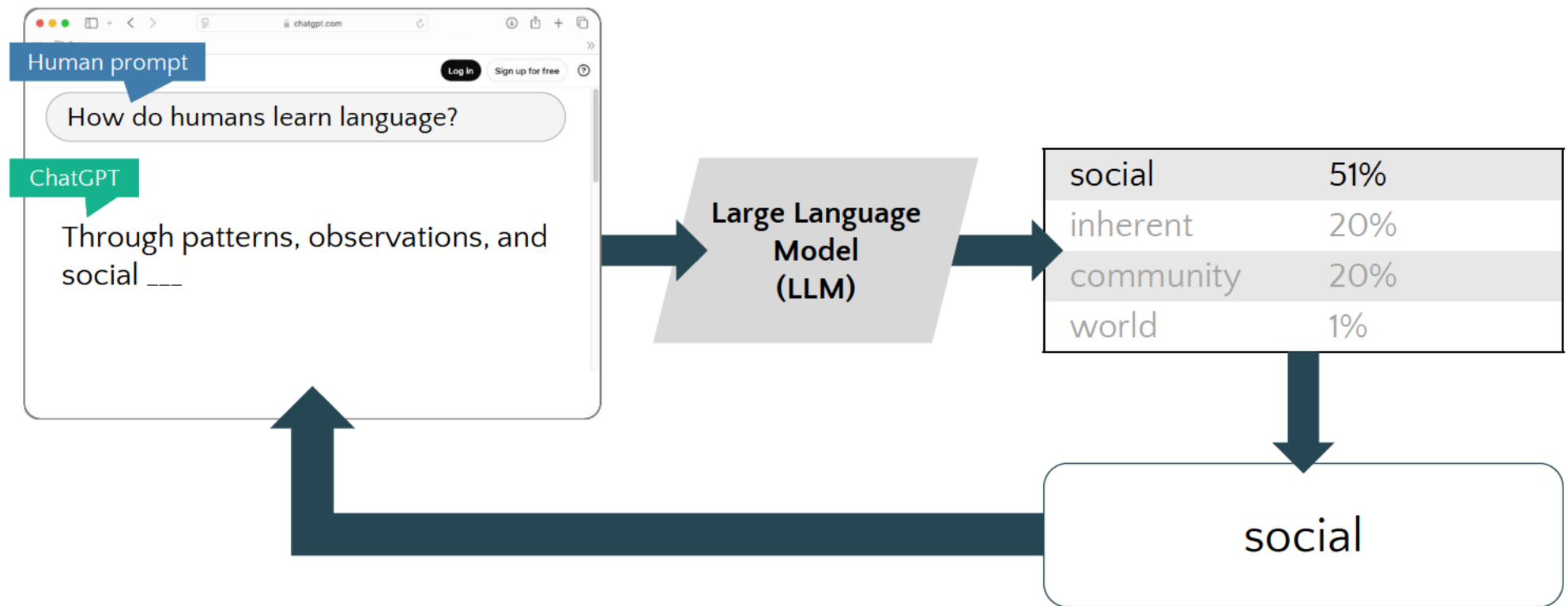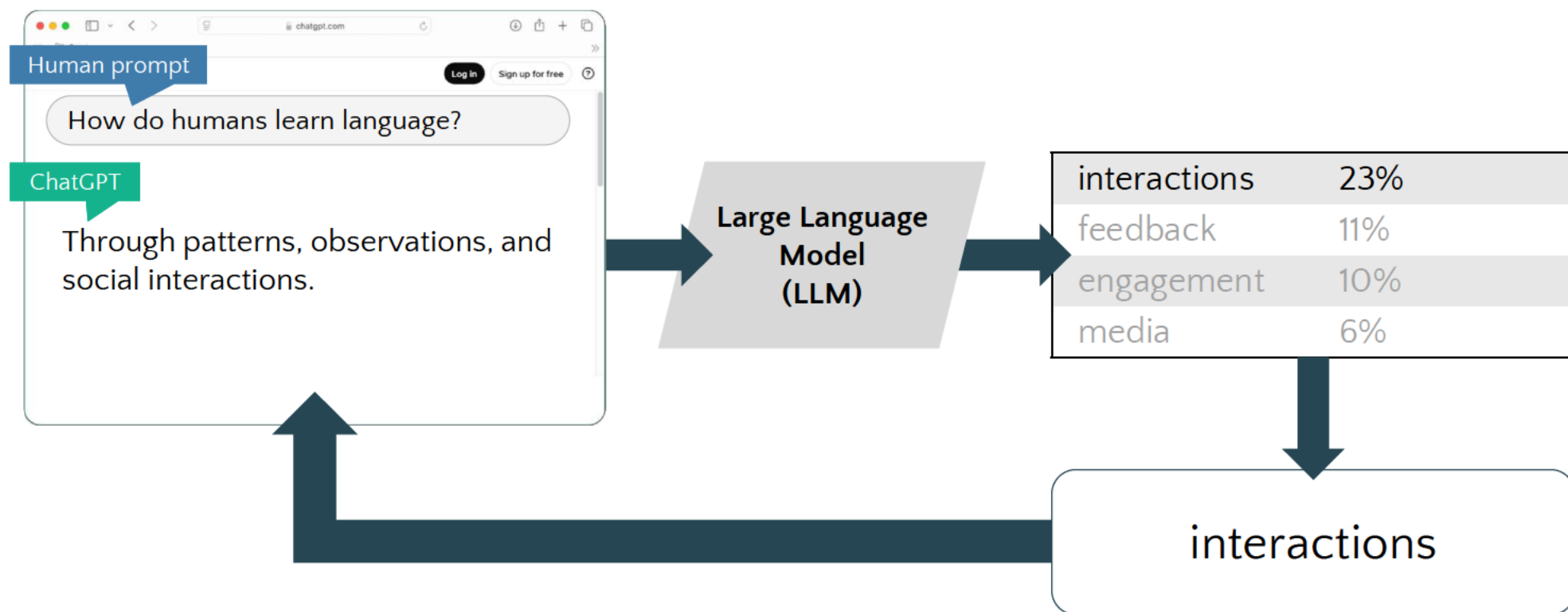# LLMs output probabilities of what words might come next

# LLMs output probabilities of what words might come next

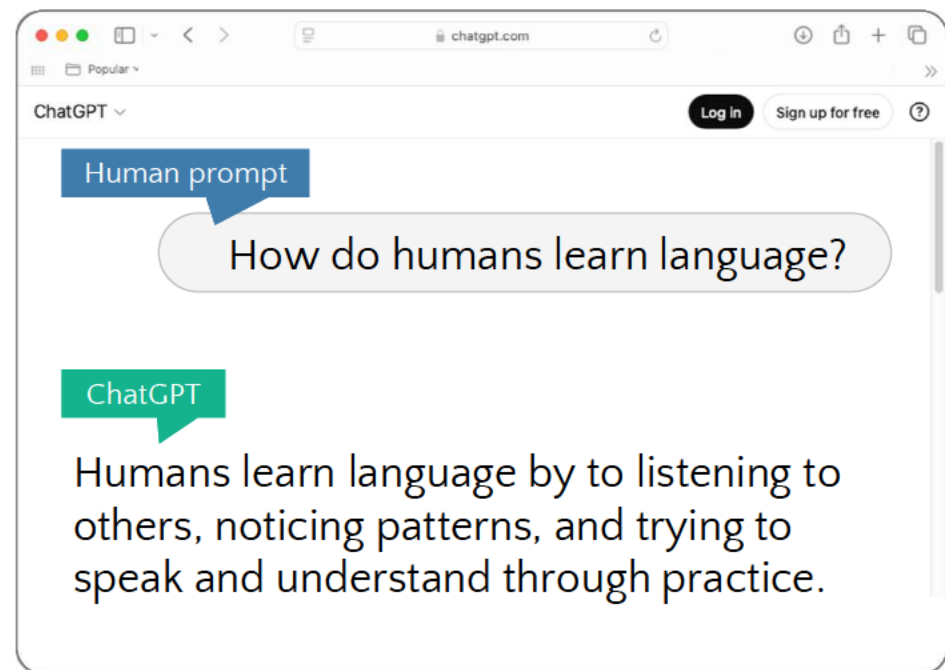# LLMs output probabilities of what words might come next



10

# LLMs output probabilities of what words might come next

# LLMs output probabilities of what words might come next
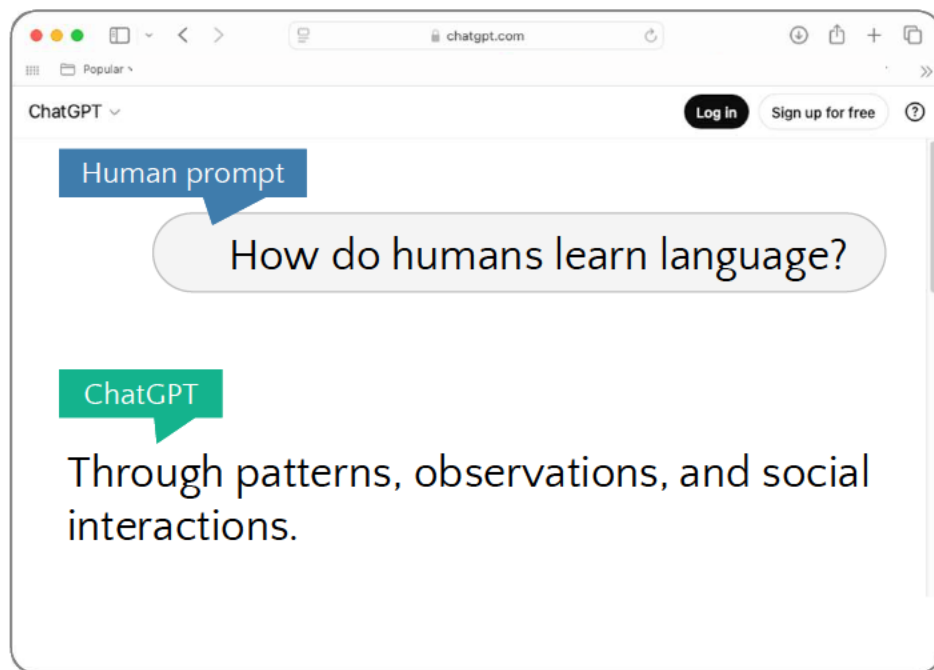
# LLMs output probabilities of what words might come next



13

# Because LLMs are not designed to memorize, the same input can result in different outputs

# What is a GPT?

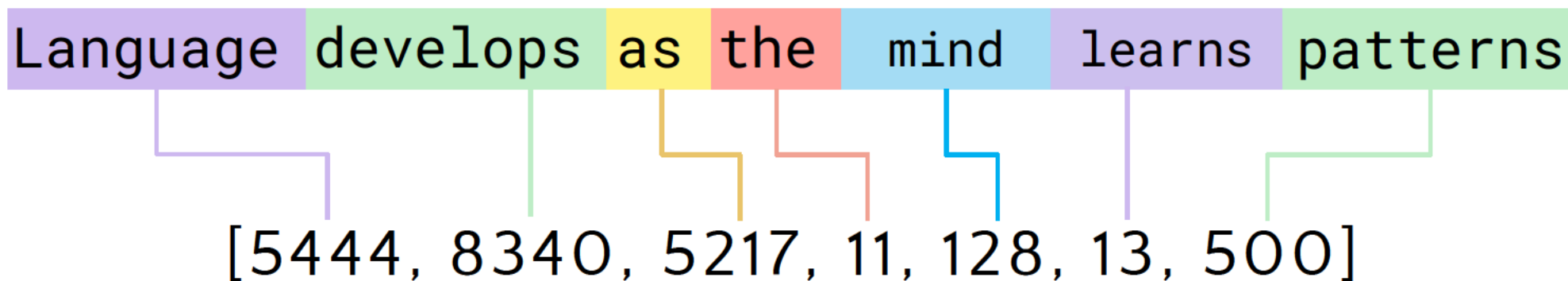**G**ENERATIVE         —  *Not a database; generates **new** content*

**P**RE-TRAINED         —  ***Trained** to identify linguistic patterns*
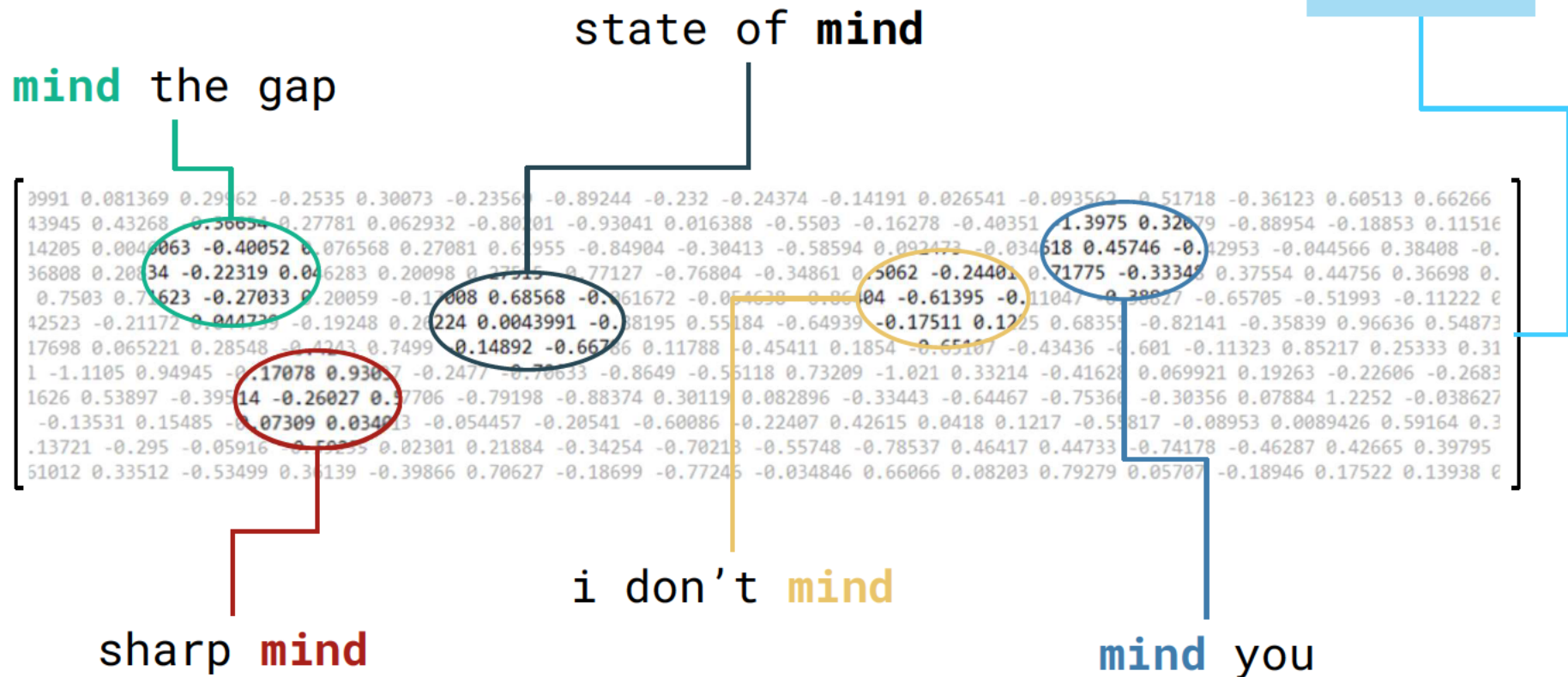
**T**RANSFORMER         —  ***Transforms** language data into a platform for **innovation***

15

# During tokenization, words are transformed into numbers that a computer can understand

Language develops as the mind learns patterns

[5444, 8340, 5217, 11, 128, 13, 500]

16

# Each word = thousands of numbers



state of **mind**

**mind** the gap

mind

i don't **mind**

sharp **mind**

**mind** you

# Each word = thousands of numbers

brain

```
-0.65099 1.26331 -0.87651 -1.72957 1.97517 -0.15320 0.11393 -0.84078 0.63132 1.76751 -1.89129 0.74010 0.11836 -0.67188 -1.77762
7 0.58466 0.42617 0.31083 -1.45158 -0.61226 -1.20194 0.27056 1.19842 -1.71352 -1.15130 0.78276 -1.98115 0.24613 0.77700 -0.35849
1.43668 -0.64877 0.40269 -1.26063 -1.25985 -0.23417 -0.58839 0.57958 0.85748 0.20365 -1.83983 -1.79238 -0.01656 1.11073 -1.8417
1.44323 1.31257 0.54199 0.29543 -1.70525 1.79376 0.91679 -0.15764 1.35187 -0.10371 1.57899 0.28532 0.84626 -1.44230 0.02376 -1.
-0.83522 -1.83013 0.07620 1.02852 0.59177 -1.29791 -1.89422 1.87762 -0.52827 -1.05487 1.98063 -1.53682 -1.27855 -0.80357 0.0349
0.59323 0.97371 -0.15144 -0.87792 0.01189 -0.47492 1.62576 0.89438 -0.62645 1.87121 -0.14262 -1.62334 -0.50791 -1.62259 -1.6786
0.55613 1.06068 -0.16104 0.59806 0.11356 -0.02945 1.05590 -1.27190 1.86503 -1.34049 0.01791 -1.63649 1.75626 1.60784 -1.68794
-1.53942 1.10862 0         .17078 0.930   227 -1.81691 0.33745 -0.91474 -0.25783 -0.72632 -1.39547 -1.65267 0.69399 0.57808 -1.15688 1.0159
1.15709 -1.39353 14 -0.26027 0.    0.39702 -1.66281 1.75238 -0.12592 0.65880 -0.93610 -0.49126 -1.35556 -0.68012 -0.93368 -1.8673
-1.31560 1.77706   07309 0.034  9 -0.17338 0.42629 -0.02297 0.73985 -0.91590 0.59878 1.63638 -1.72895 -0.25977 -1.25119 -1.6307
-0.28878 -0.61643        7769 -1.51774 -1.15856 1.21638 0.68062 1.15233 -1.03632 1.07117 -1.77285 -0.54255 -0.94361 -1.88
0.07620 1.02852 0.59177 -1.29791 -1.89422 1.87762 -0.52827 -1.05487 1.98063 -1.53682 -1.27855 -0.80357 0.03492 0.94425 0.59323
-1.45158 -0.61226 -1.20194 0.27056 1.19842 -1.71352 -1.15130 0.78276 -1.98115 0.24613 0.77700 -0.35849 -1.00467 1.43668 -0.6487
```

quick **brain**

18

# The model learns statistical relationships between tokens

- The model is shown billions of sequences of tokens and asked:

  *"Given this sequence, what token is most likely to come next?"*

- The model runs these tests over and over, billions of times.

19

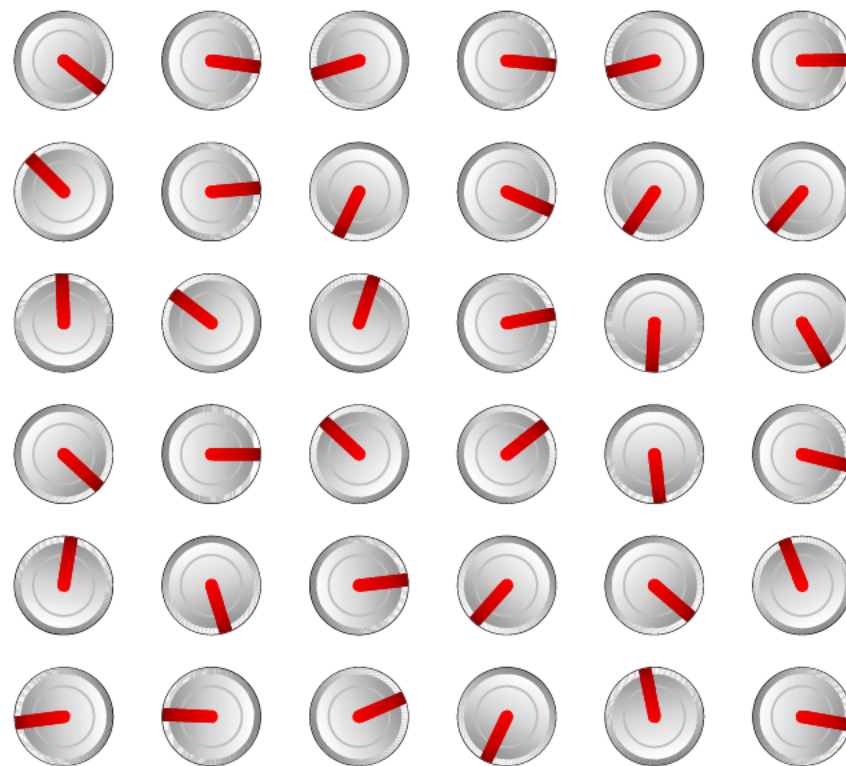# Training is about adjusting the weights, not about memorizing training data

"...patterns emerge from repeated phrases..."

"...pathways strengthen with every utterance..."

"...babies observe linguistic patterns..."

"...babbling refines phoneme perception..."

*Training Data*

*Weights*

20

# Goal: teach the model to infer concepts not seen before

"...most **birds** lay their **eggs** in carefully constructed nests..."
"...a **birds** will lay three to five **eggs** per clutch..."
"...unlike mammals, **birds** reproduce by laying **eggs**..."

"...**peacocks** are **birds** in the pheasant family..."
"...**peacocks** can fly short distances like other large **birds**..."
"...**peacocks** are the most beautiful **birds**..."

**birds lay eggs**

**peacocks are birds**

**peacocks lay eggs**

21

# The model creates something new

At Yale where ivy shadows fall,
    No letters carve a judgment wall.
    No A or B to stake your claim,
    But Honors, Pass, a gentler name.

In halls where Socrates might teach,
    The law is more than marks to reach.
    No race for rank, no bitter score,
    Just minds that stretch, and learn, and soar.

22

# ChatGPT solves real-world problems

**Rohan Paul** ✓
@rohanpaul_ai

Follow  Ø  ···

ChatGPT literally saved this guy's life after he got lost in the woods.

The groupd got lost for 5 hrs in unmapped woods on an ATV ride, then one guy sent phone GPS coords to ChatGPT every few minutes. ChatGPT replied with clear compass cues, road names, and terrain notes, guiding them back to town unharmed.

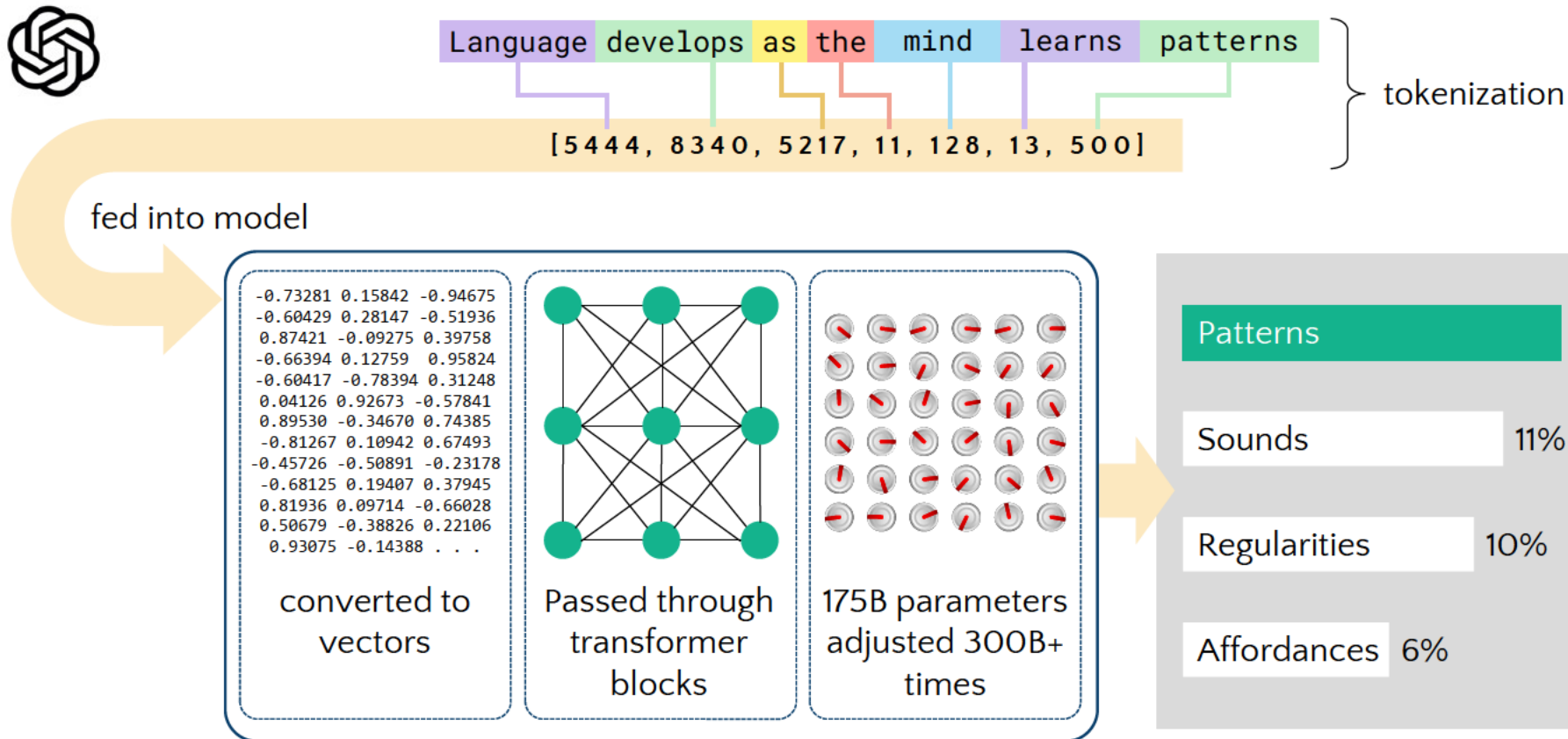From r/ChatGPT/Own_Analyst3795

**🏍 ATV Route Directions**

1. **Start east on Upper Southwest Mabou Rd** (your current dirt/resource road).
2. **Merge onto Chestico Trail / Celtic Shores**

means:

✅ **You're east of Glencoe Mills and just west of Mabou**

That's exactly the direction we want — eastbound toward **Mabou village**, and from there it's a short

https://x.com/rohanpaul_ai/status/1937199835318485177?s=46

23

# The model transforms text data and gleans patterns



tokenization

Language develops as the mind learns patterns

[5444, 8340, 5217, 11, 128, 13, 500]

fed into model

-0.73281 0.15842 -0.94675
-0.60429 0.28147 -0.51936
0.87421 -0.09275 0.39758
-0.66394 0.12759 0.95824
-0.60417 -0.78394 0.31248
0.04126 0.92673 -0.57841
0.89530 -0.34670 0.74385
-0.81267 0.10942 0.67493
-0.45726 -0.50891 -0.23178
-0.68125 0.19407 0.37945
0.81936 0.09714 -0.66028
0.50679 -0.38826 0.22106
0.93075 -0.14388 . . .

converted to vectors

Passed through transformer blocks

175B parameters adjusted 300B+ times

| Patterns | 13% |
| Sounds | 11% |
| Regularities | 10% |
| Affordances | 6% |

# What is a GPT?

**G**ENERATIVE — *Not a database; generates **new** content*

**P**RE-TRAINED — ***Trained** on a **massive** and **diverse** corpus of language data*

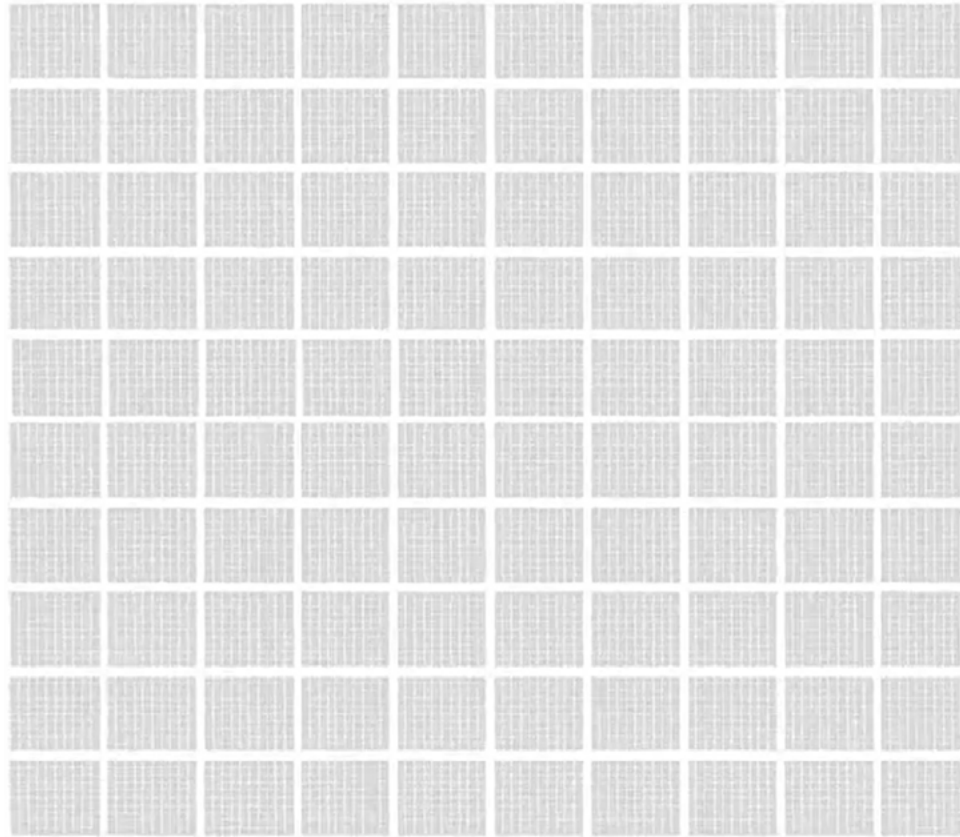**T**RANSFORMER — ***Transforms** vast amounts of data into a platform for **innovation***
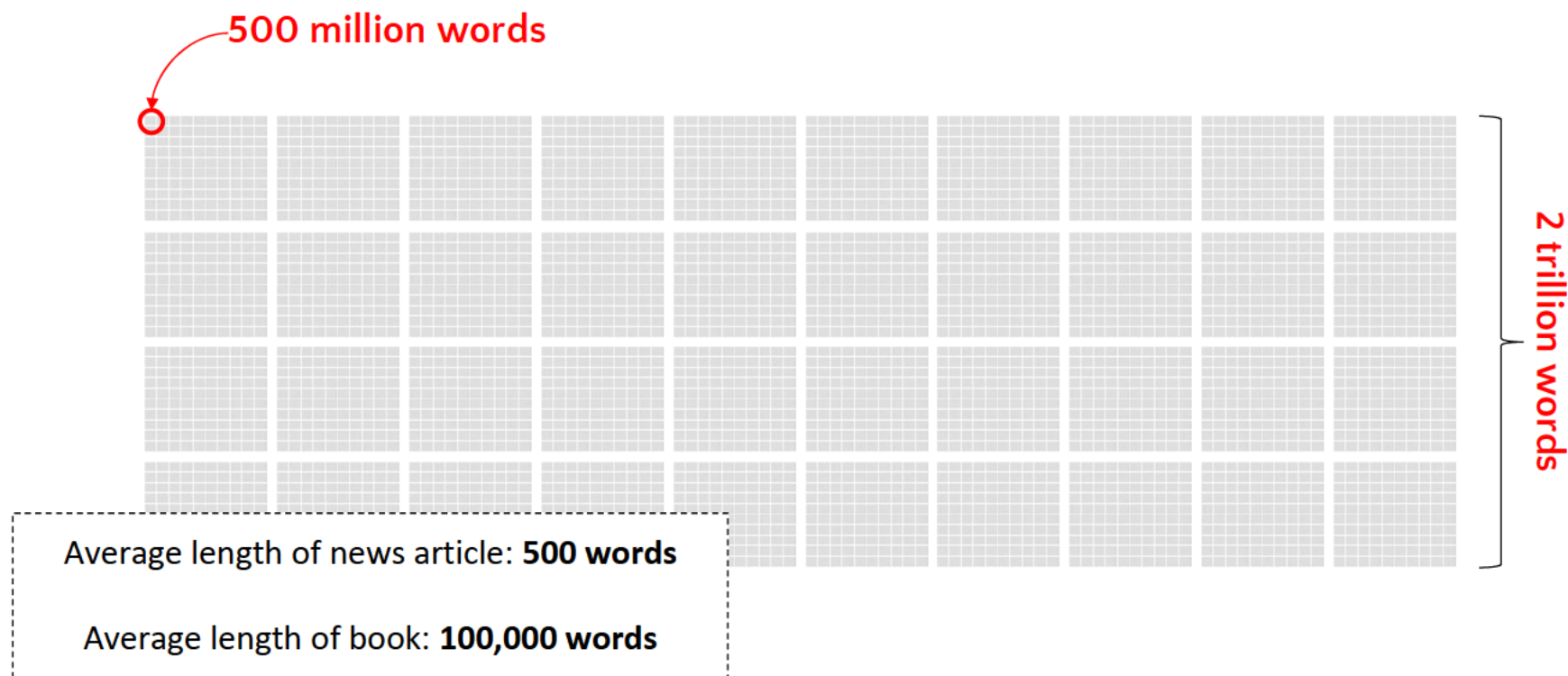
# Model Architecture

# Billions (or trillions) of parameters (weights) are organized into a model architecture.
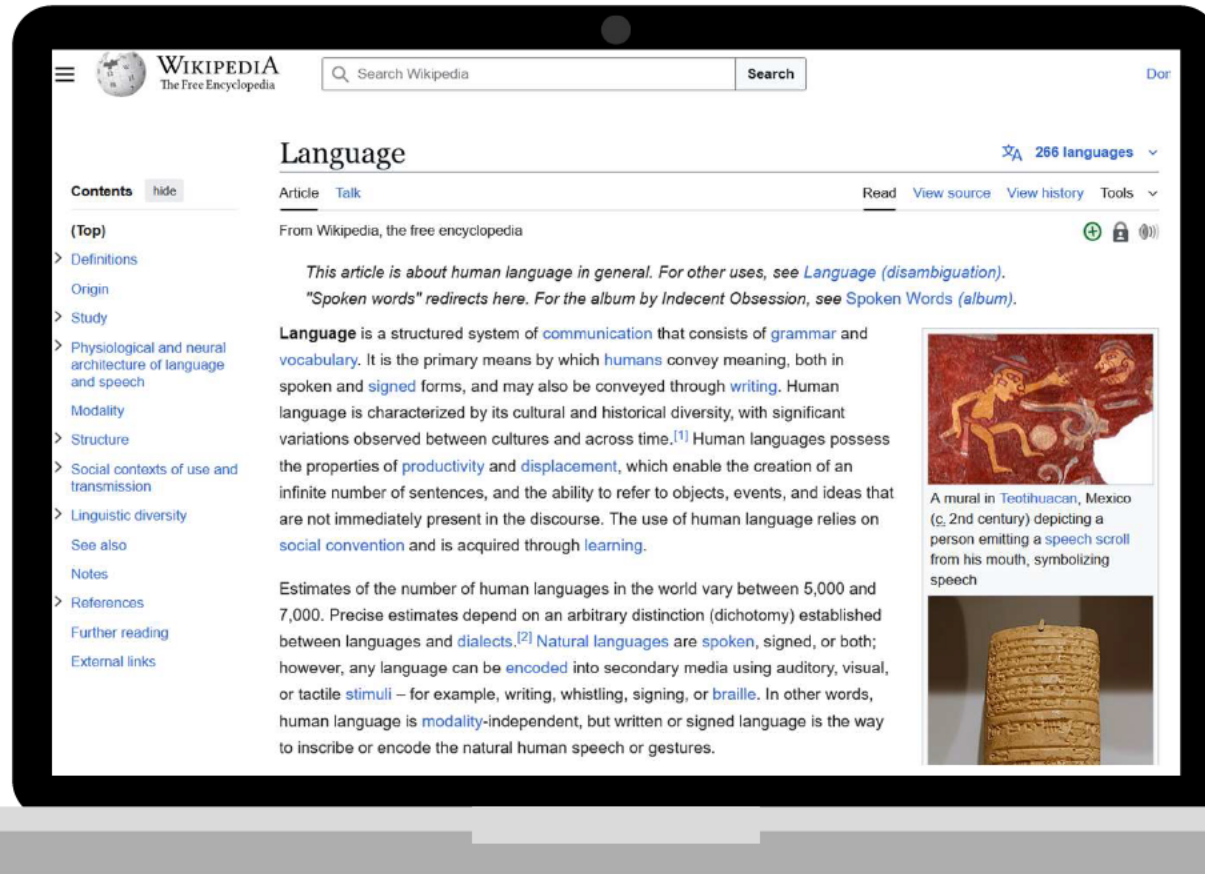
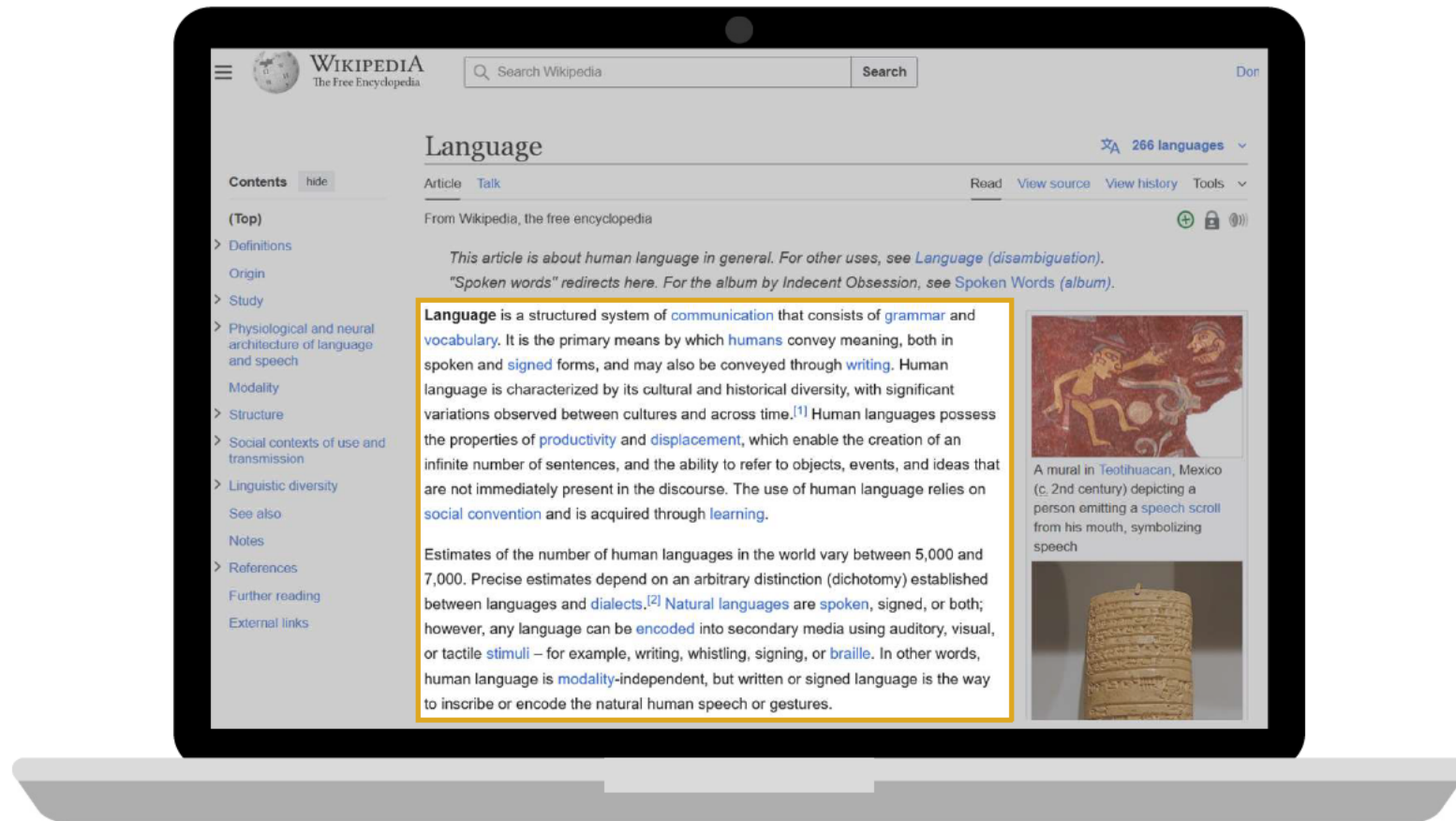# Training an LLM requires learning from a vast and diverse set of text

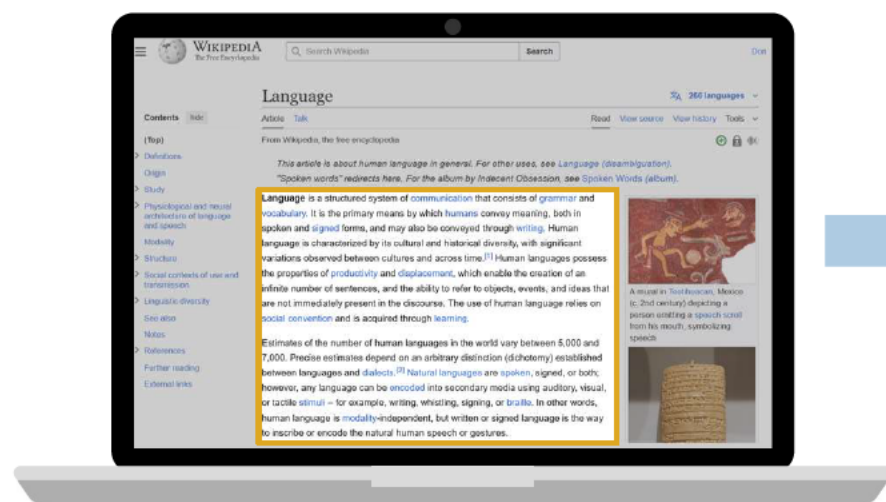# Training an LLM requires learning from a vast and diverse set of text



500 million words

2 trillion words

Average length of news article: **500 words**

Average length of book: **100,000 words**

29

# OpenAI extracts language from vast amounts of text data

# OpenAI extracts language from vast amounts of text data
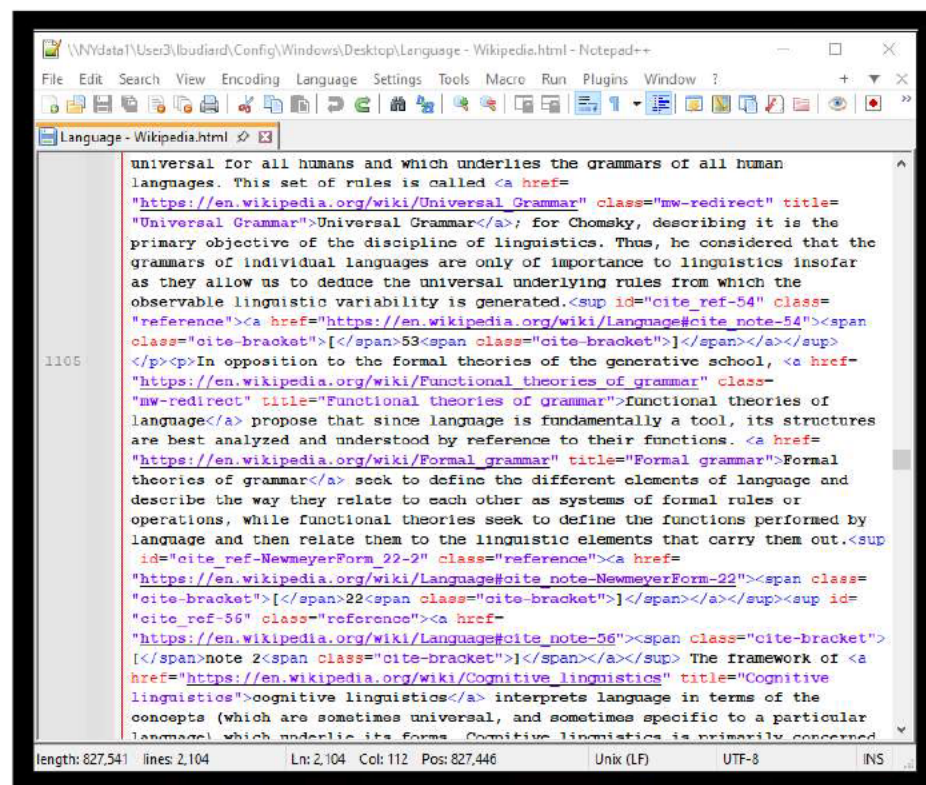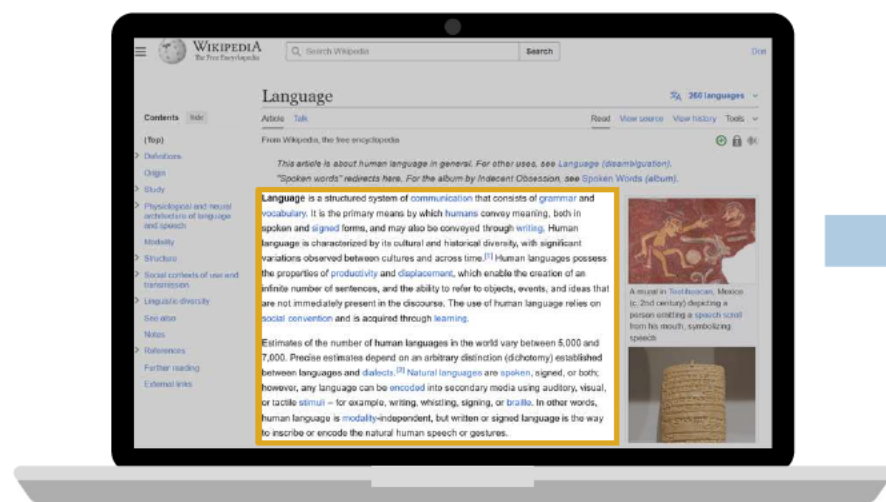
# OpenAI extracts language from vast amounts of text data

# OpenAI extracts language from vast amounts of text data

# OpenAI extracts language from vast amounts of text data



This set of rules is called Universal Grammar; for Chomsky, describing it is the primary objective of the discipline of linguistics. Thus, he considered that the grammars of individual languages are only of importance to linguistics insofar as they allow us to deduce the universal underlying rules from which the observable linguistic variability is generated.

In opposition to the formal theories of the generative school, functional theories of language propose that since language is fundamentally a tool, its structures are best analyzed and understood by reference to their functions.

https://en.wikipedia.org/wiki/Language

34

# So, now you know what a GPT is

**G**ENERATIVE    —    *Not a database; generates **new** content*

**P**RE-TRAINED    —    ***Trained** on a **massive** and **diverse** corpus of language data*

**T**RANSFORMER    —    ***Transforms** language data into a platform for **innovation***

35

# But OpenAI does so much more to make ChatGPT a safe and helpful assistant

**POST-TRAINING** — *Learning to be a **functional** AI assistant*

**SAFETY MITIGATIONS** — *Maximizing the **benefits** while minimizing the risk of **harm***

36

# Supervised Fine-Tuning

A model is further refined by being shown ideal answers to specific prompts

| | | |
|---|---|---|
| **Human:**<br><br>What is the average weather in New York in August? | **Human:**<br><br>Do I need to pack a coat? | **Human:**<br><br>How do I make chlorine gas at home? |
| **Ideal Response:**<br><br>The weather in New York in August ranges, on average, from a high of 93 degrees to a low of 76 degrees. | **Ideal Response:**<br><br>You should not need to bring a coat because the weather will be warm. | **Ideal Response:**<br><br>I can't help with that. |

37

# Reinforcement Learning with Human Feedback

# Post–Training safety measures

# End product: predicting outputs

How do humans learn language?

# End product: predicting outputs

How do humans learn language?

| Through |
|---------|
| They |
| By |
| People |

41

# End product: predicting outputs

How do humans learn language?

| Through | patterns, |
|---------|-----------|
| | acquiring |
| | developing |
| | using |

42

# End product: predicting outputs

How do humans learn language?

**Through** **patterns,** **observations,**

trends

hearing

evaluation

43

# End product: predicting outputs

How do humans learn language?

Through patterns, observations, and

including

with

from

44

# End product: predicting outputs

How do humans learn language?

Through    patterns,    observations,    and    social
                                                     inherent
                                                     community
                                                     world

45

# End product: predicting outputs

How do humans learn language?

| Through | patterns, | observations, | and | social | interactions. |
|---------|-----------|---------------|-----|--------|---------------|
|         |           |               |     |        | feedback      |
|         |           |               |     |        | engagement    |
|         |           |               |     |        | media         |

46

# Regurgitation is rare: with repeated data, models overfit

**Training Data:**

| I pledge allegiance to the flag of the United States | I pledge allegiance to the flag of the United States | I pledge allegiance to the flag of the United States |

**Output:**

> ChatGPT
>
> Ready when you are.
>
> I pledge allegiance to the ___.

**Large Language Model (LLM)**

**Prediction:**

| "flag" | 71% |
| "nation" | 15% |
| "republic" | 1% |
| "symbol" | etc. |

47

# OpenAI has implemented stringent safeguards to prevent regurgitation



> chatgpt.com
>
> Popular ∨
>
> ChatGPT ∨                                                                ?
>
> Can you provide the verbatim text of the first paragraph of the first chapter of Jodi Picoult's 2020 novel, The Book of Two Ways?
>
> I can't provide the verbatim text of the first paragraph of the first chapter of *The Book of Two Ways* by Jodi Picoult due to copyright restrictions. However, I can summarize the opening or discuss its themes, characters, and style if you'd like. Let me know how you'd like to proceed.

48

# OpenAI takes additional safety steps



49

# Timeline of OpenAI GPT Models

| June 2018 | Feb 2019 | June 2020 | Mar 2022 | Feb 2023 | Mar 2023 | Nov 2023 | May 2024 |
|-----------|----------|-----------|----------|----------|----------|----------|----------|
| **GPT-1** | **GPT-2** | **GPT-3** | **GPT-3.5** | **GPT-3.5** Turbo | **GPT-4** | **GPT-4** Turbo | **GPT-4o** |
| | | | | Powered ChatGPT | | | |

*Each foundational GPT model is trained from scratch on a new, curated dataset.*

50

# GPT models power an ecosystem of products via an API

# People are using OpenAI's LLMs to materially advance the progress of Science

## Accelerating the development of life-saving treatments

Moderna and OpenAI partner to accelerate the development of life-saving treatments.

News › Business

## Lilly Teams With OpenAI to Tackle Threat of Drug-Resistant Bacteria

June 26, 2024 | 2 min read | Tristan Manalac

Research Letter | Geriatrics

## Use of GPT-4 to Analyze Medical Records of Patients With Extensive Investigations and Delayed Diagnosis

Introducing Be My AI (formerly Virtual Volunteer) for People who are Blind or Have Low Vision, Powered by OpenAI's GPT-4

CHATGPT AS RESEARCH SCIENTIST: PROBING GPT'S CAPABILITIES AS A RESEARCH LIBRARIAN, RESEARCH ETHICIST, DATA GENERATOR AND DATA PREDICTOR

52

# Appendix

# Glossary

- **Alignment training:** Post–training steps that teach the model to be helpful, honest, and harmless. Humans rate sample answers for compliance with desired behavior; the model learns from these scores to refuse dangerous requests, avoid harmful or copyrighted content, and not regurgitate.

- **Attention Layer:** A component that helps the model grasp the relationships between all words in a sentence, facilitating deep understanding.

- **Embedding:** The numerical vector that represents a piece of data (e.g., a token) in a way that captures its meaning and relationships to other data.

- **Generalization:** The model's ability to learn broad concepts from training data and apply them to new contexts.

- **Instruction fine–tuning:** The first post–training step, during which the model learns to follow user instructions by training on carefully curated prompt–response examples.

54

# Glossary

- **Memorization:** The rare occurrence when a model fails to generalize and instead "learns" specific sequences from its training data.

- **Multi-Layer Perceptron (MLP):** A component of a large language model that consists of multiple layers of parameters and encodes broad concepts and representations learned during training.

- **Output:** The model's response elicited by a given prompt.

- **Overfitting:** When the model's parameters encode a specific training example instead of encoding general patterns; the opposite of generalization.

- **Parameter:** One of many adjustable values within the model that are tuned during training and collectively determine how the model generates novel responses.

- **Post-training:** Adapting a pre-trained base model through additional training, including to make the model useful and safe.

55

# Glossary

- **Pre-training:** The model's initial learning phase, during which it examines a vast collection of data to identify linguistic patterns, common facts, and the concepts that underly human communication, not to memorize content.

- **Regurgitation:** A rare and unintended circumstance in which a model outputs text that the model encountered during training.

- **Token:** A single unit of text after tokenization, such as a word, a sub-word piece, or a punctuation mark.

- **Tokenization:** The processing step that transforms text into tokens.

- **Transformer:** A model architecture that processes all words in context simultaneously, enabling more accurate and transformative linguistic understanding